

## CLAIMS

What is claimed is:

1 1: A method for managing distribution of messages for changing the state of shared data in a  
2 computer system having a main memory, a memory management system, a plurality of processors,  
3 each processor having an associated cache, and employing a directory-based cache coherency  
4 comprising the method of:

5 grouping the plurality of processors into a plurality of clusters;  
6 tracking copies of shared data sent to processors in the clusters;  
7 receiving an exclusive request from a processor requesting permission to modify a shared  
8 copy of the data;

9 generating invalidate messages requesting that other processors sharing the same data  
10 invalidate that data;

11 sending the invalidate messages only to clusters actually containing processors that have a  
12 shared copy of the data in the associated cache; and

13 broadcasting the invalidate message to each processor in the cluster.

1 2. The method of claim 1, wherein the invalidate message is sent to one master processor in a  
2 cluster, and further comprising:

3 the master processor distributing the invalidate message to one or more slave processors  
4 and waiting for an acknowledgement from said one or more processors;

5 if said one or more slave processors are configured to do so, distributing the invalidate  
6 message to one or more other slave processors, if any exist, and waiting for an acknowledgement  
7 from said other slave processors;

8           a slave processor which does not distribute the invalidate message to any other processor  
9        replying with an acknowledgement to the processor from which the invalidate message was  
10      received; and

11           upon receiving acknowledgements from all processors to which the invalidate messages  
12      were sent, a slave processor replying with an acknowledgement to the processor from which the  
13      invalidate message was received;

14           wherein upon receiving an invalidate message, the processor invalidating a local copy of  
15      the shared data, if it exists, and wherein upon receiving acknowledgements from all slave  
16      processors to which the invalidate messages were sent, the master processor sending an invalidate  
17      acknowledgment message to the processor that originally requested the exclusive rights to the  
18      shared data.

1    3.      The method of claim 2, wherein:

2           the slave processors to which the master processor distributes the invalidate message are  
3      determined by data registers associated with the master processor; and

4           any other slave processors to which the slave processors distribute the invalidate message  
5      are determined by data registers associated with each slave processor;

6           wherein data registers exist and may be unique for each processor entry port.

1    4.      The method of claim 2, wherein:

2           tracking of the shared copies of the data sent to the clusters is performed by setting a bit in  
3      a data register with at least as many bit positions as there are clusters;

4           wherein each cluster is associated with one bit position in the data register.

1       5.    The method of claim 4, wherein sending the invalidate messages only to one master  
2   processor in a cluster actually containing processors that have a shared copy of the data in the  
3   associated cache further comprises the steps of:

4           selecting only the bit positions containing a set bit;  
5           cross referencing the bit positions with cluster numbers;  
6           cross referencing cluster numbers with an actual processor identification; and  
7           delivering the invalidate message to the processor associated with the processor  
8   identification.

1       6.    The method of claim 1, further comprising:  
2           distributing the main memory among and coupled to each of the plurality of processors and  
3   each processor comprising a directory controller for the main memory coupled to that processor;  
4           the directory controller managing the main memory location for the shared data and  
5   tracking the copies of shared data sent to processors in the clusters;  
6           the processor requesting exclusive ownership of the shared data delivering the request to  
7   the directory controller; and  
8           the directory controller sending the invalidate messages to master processors in clusters  
9   actually containing processors that have a shared copy of the data.

1       7.    The method of claim 1, further comprising:

2       upon receiving a request from a processor requesting permission to modify a shared copy  
3    of the data, sending a response to the requesting processor indicating the number of additional  
4    shared copies of the data;

5       changing the state of the shared data by the requesting processor from shared to exclusive;  
6    and

7       waiting to modify the exclusive data until acknowledgements arrive from the clusters  
8    actually containing processors that have a shared copy of the data in the associated cache.

1   8.   The method of claim 7, wherein:

2       when the processor requesting exclusive ownership of the shared data, the directory  
3    controller, and shared copies of the data exist within the same cluster, the directory node assumes  
4    the position of master node and broadcasts the invalidate message to all the processors in the  
5    cluster.

1   9.   A multiprocessor system, comprising:

2       a main memory configured to store data;

3       a plurality of processors, each processor coupled to at least one memory cache;

4       a memory directory controller employing directory-based cache coherence;

5       at least one input/output device coupled to at least one processor;

6       a share mask comprising a data register for tracking shared copies of data blocks that are

7    distributed from the main memory to one or more cache locations; and

8       a PID-SHIFT register which stores configuration settings to determine which one of several  
9    shared data invalidation schemes shall be implemented;

10 wherein when the PID-SHIFT register contains a value of zero, the data bits in the share  
11 mask data register correspond to one of the plurality of processors and wherein when the PID-  
12 SHIFT register contains a nonzero value, the data bits in the share mask data register correspond to  
13 a cluster of processors, each cluster comprising more than one of the plurality of processor.

1 10. The system of claim 9 wherein:

2 if the value in the PID-SHIFT register is zero, the directory controller sets the bit in the  
3 share mask corresponding to the processor to which a shared copy of a data block is distributed;  
4 and

5 wherein if the value in the PID-SHIFT register is nonzero, the directory controller sets the  
6 bit in the share mask corresponding to the cluster containing a processor to which a shared copy of  
7 a data block is distributed.

1 11. The system of claim 10 wherein:

2 the nonzero value in the PID-SHIFT register determines the number of processors in each  
3 cluster.

1 12. The system of claim 9 wherein:

2 when more than one shared copy of a data block exists outside of the main memory; and  
3 wherein in response to a request from a requesting processor for exclusive write access to  
4 one of the shared copies of the data block; and

5           wherein when the value in the PID-SHIFT register is zero, the directory controller transmits  
6   an invalidate message only to those processors whose corresponding bits in the share mask are set,  
7   except the requesting processor; and

8           wherein when the value in the PID-SHIFT register is nonzero, the directory controller  
9   transmits an invalidate message only to those clusters whose corresponding bits in the share mask  
10   are set.

1   13.   The system of claim 12 wherein the cluster further comprises:

2           a master processor to which the invalidate message directed toward the cluster are  
3   delivered; and

4           one or more slave processors, each of which receive an invalidate message that is generated  
5   by the master processor.

1   14.   The system of claim 13 further comprising:

2           a processor router table that includes cross reference information which correlates master  
3   processor identification with cluster numbers.

1   15.   The system of claim 13 further comprising:

2           configuration registers associated with each port of a processor in a cluster which  
3   determine the path by which the invalidate message is broadcast within a cluster.

1   16.   A multiprocessor system, comprising:

2           a memory;

3       multiple computer processor nodes, each with an associated memory cache; and  
4       a memory controller employing a directory-based cache coherency employing shared  
5       memory invalidation method, wherein:

6       the nodes are grouped into clusters;

7       the memory controller distributes memory blocks from the memory to the various cache  
8       locations at the request of the associated nodes;

9       upon receiving a request for exclusive ownership of one of the shared memory blocks, the  
10      memory controller distributes invalidate messages via direct point to point transmission to only

11      those clusters containing nodes that share a block of data in the associated cache; and

12      wherein when the invalidate message is received by a cluster, an invalidate message is  
13      broadcast to all nodes in the cluster.

1       17.   The system of claim 16 further comprising:

2       a share mask data register with as many bit locations as there are clusters;

3       a router lookup table with cross reference information correlating bit locations in the share  
4       mask to one master nodes in each cluster;

5       wherein the memory controller determines to which cluster to send the invalidate message  
6       according to bits set in the share mask and sends the invalidate message to the router which then  
7       forwards the invalidate message to the node whose identification corresponds to the cluster number  
8       as indicated in the router table.

1       18.   The system of claim 17 each node further comprising:

2           router control and status registers for each input port of the node which configure the  
3   node's broadcast forwarding scheme wherein the forwarding scheme determines to which, if any,  
4   nodes the node shall forward a broadcast invalidate message when a broadcast invalidate message  
5   is received at a given port.

1   19.    The system of claim 18 wherein:  
2           the router control and status registers are comprised of bit locations corresponding to each  
3   output port of the node; and  
4           wherein if a bit location contains a set bit, the invalidate message is forwarded to the output  
5   port corresponding to that bit location; and  
6           wherein if a bit location does not contain a set bit, the invalidate message is not forwarded  
7   to the output port corresponding to that bit location.

1   20.    The system of claim 19 wherein the processors in a cluster invalidate shared data, if it  
2   exists, and generate and forward acknowledgments in reverse direction but along the same path  
3   followed by the invalidate messages.